

01

Into a Parallel New World

Yangdong Steve Deng (邓仰东)
Tsinghua University, Beijing



Parallel World



Parallel Brain

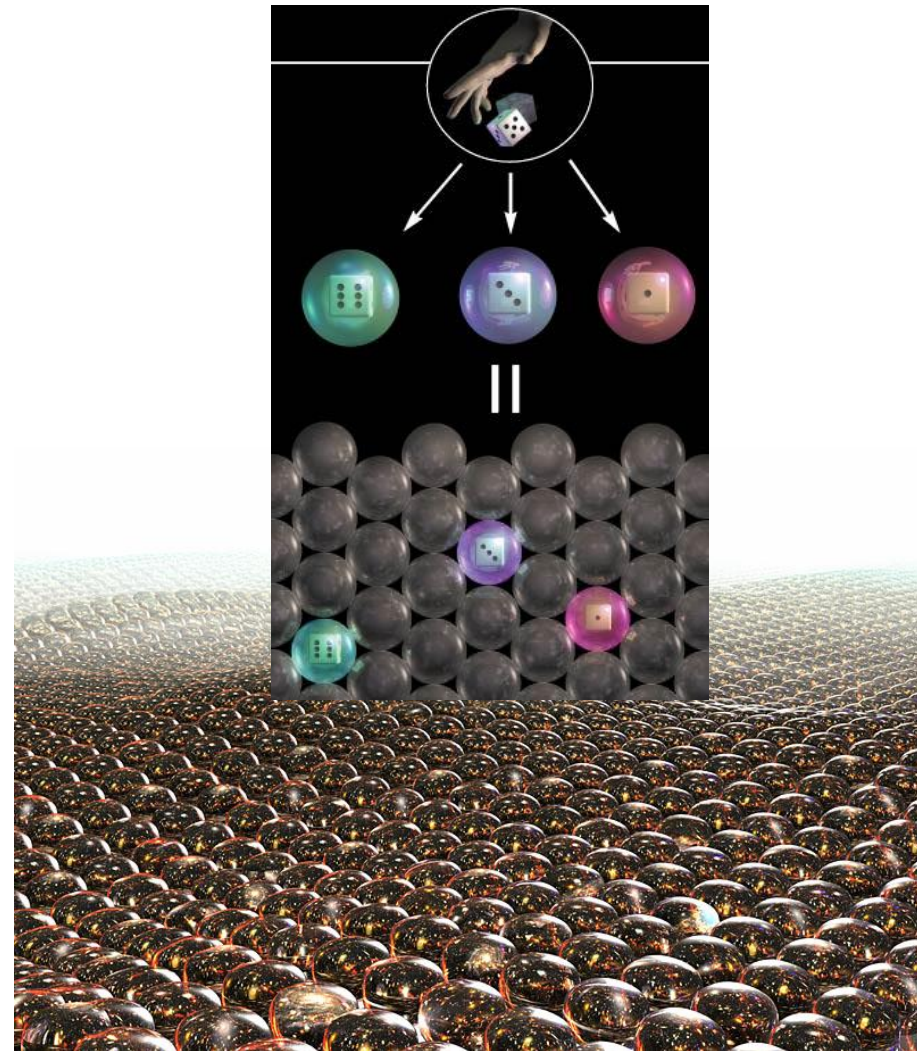


blue neurons / neu / January 2010.

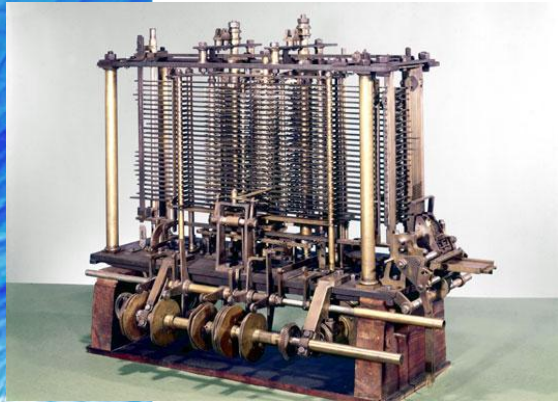
Parallel Universe



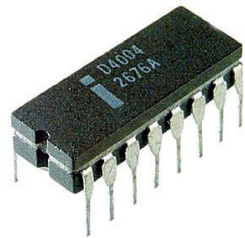
Schrodinger's cat



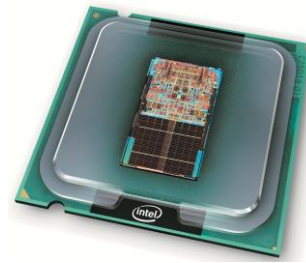
The Evolution to Parallel Processors



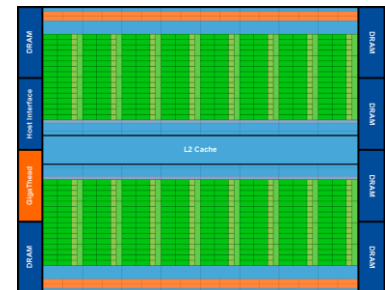
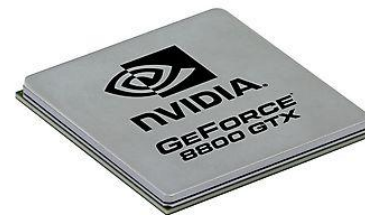
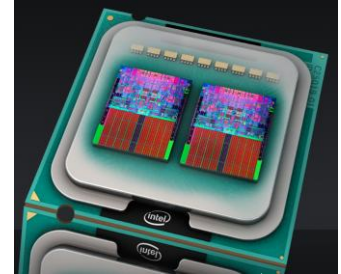
**Analytical Engine
(1871)**



**Single Core
(1971)**



**Multi/Many Core
(2006)**



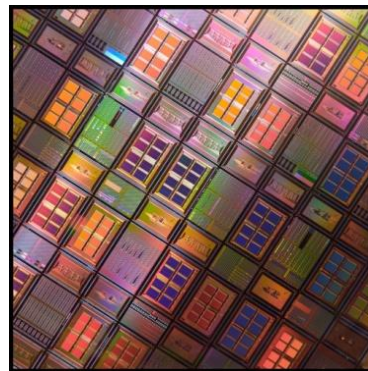
Outline



清华大学

Tsinghua University

1. Introduction



2. Research



3. CUDA Propagation



Map of THU (Tsinghua University)



- Line between East Gate/South Gate & Conference Hall
- Line between Conference Hall & Dining Hall
- Line between Conference Hall & Poster Hall
- I Conference Hall (Auditorium Hall of THU)(大礼堂)
- II Poster Hall (Gymnasium of THU)(综合体育馆)
- III NT09 Tutorial (Central main Bldg.)(中央主楼)
- IV Satellite Meeting Room (Eng. Physics. Bldg.)(工物馆)
- V Registraion Room (West Teaching Bldg.)(西阶)
- VI Dining Hall (Food Plaza)(饮食广场)
- VII THU Guest House (Jin Chun Yuan)(近春园)
- VIII THU Guest Huouse (Jia Suo)(甲所)

CCOE Tsinghua

- **Founded in 2009**
- **A cluster of 32 Tesla 1070s is under construction**
 - **Peak performance at 128Tflops(SP) or 16Tflops(DP)**
 - **Supporting Linpack applications on both CPU and GPU**
 - **Part of the 3rd fastest machine in China**



Research Team

- **Prof. Wenguang Chen and Prof. Guangwen Yang**
 - CS Department & High Performance Computing Center
 - High performance computing and Finance applications
- **Prof. Haixiao Gao**
 - Department of Biology
 - 3D protein structure recovery
- **Prof. Yuxiang Xing**
 - Dept. of Engineering Physics
 - CT image processing
- **Prof. Yangdong Steve Deng**
 - Institute of Microelectronics & Tsinghua-Intel Center of Advanced Mobile Computing
 - Electronic design automation and parallel computing
- **~60 students**

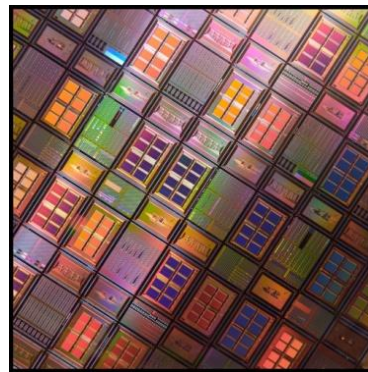
Outline



清华大学

Tsinghua University

1. Introduction



2. Research



3. CUDA Propagation

Research Outline

■ Applications and algorithms

- Logic simulation
- Packet processing
- Digital signal processing

■ Parallel microarchitecture

- A CPU+GPU integrated packet processing engine

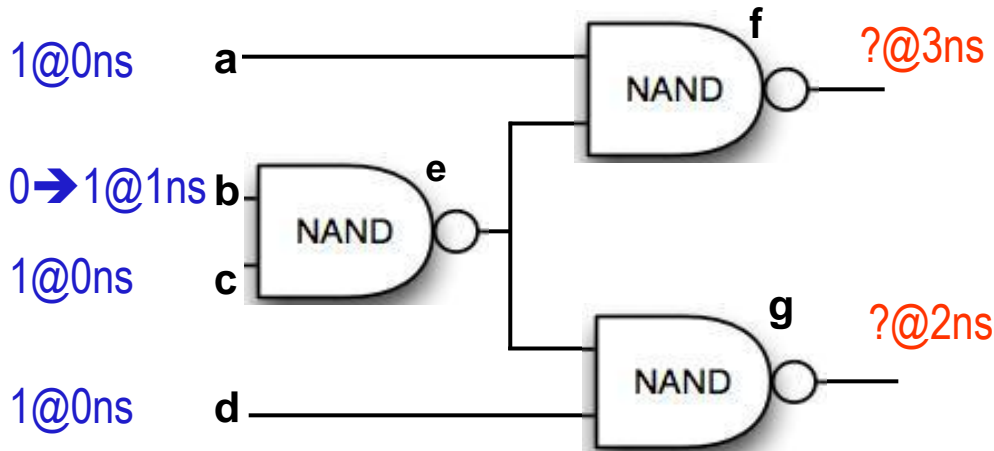
■ GPU programming

- Source-to-source GPU code generation



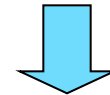
Logic Simulation

- Major method for IC design verification
 - Apply input stimuli and observe output signals
- Event driven simulation is the most widely used algorithm
 - Event: logic transition + time stamp
 - Always evaluate event with the earliest stamp

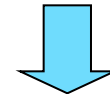


Event Queue

b: 0→1@1ns



e: 1→0@2ns

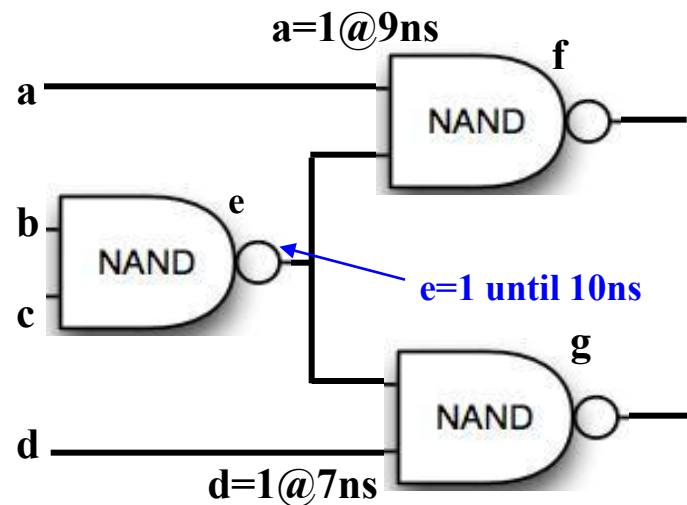


f: 0→1@3ns

g: 0→1@3ns

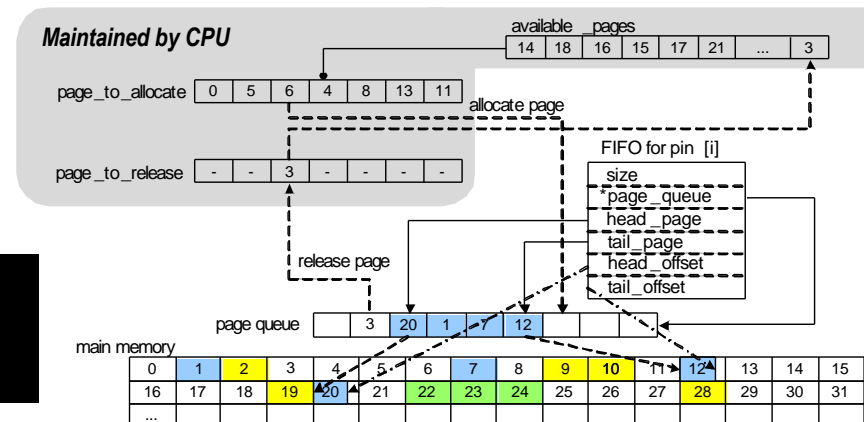
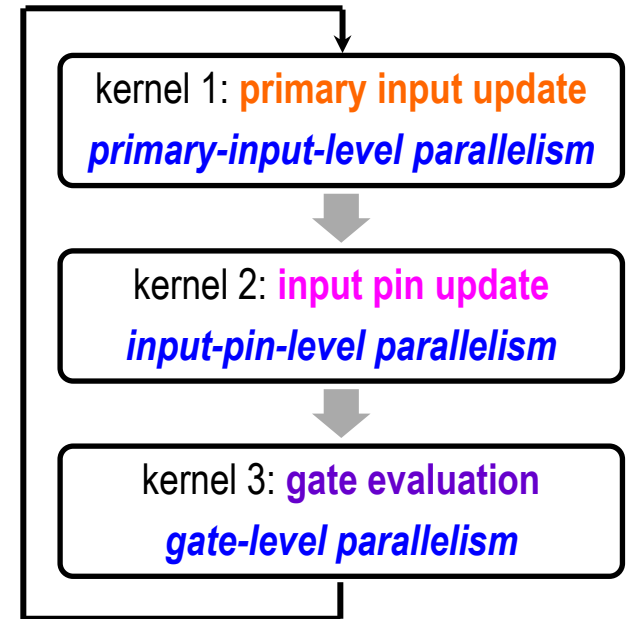
GPU Based Logic Simulation

- Simultaneously simulate events with the same time-stamp
 - Insufficient parallelism
- Chandy-Misra-Bryant (CMB) Algorithm
 - Asynchronous and conservative



Simulation Framework

- A dynamic GPU memory allocator
- World's fastest logic simulator on general purpose hardware
- 30X speed-up on average (100X for random patterns)
 - 1 month on CPU vs. 1 day on GPU



***Published on DAC 2010 and ACM Trans. TODAES 2011**

Research Outline

■ Applications and algorithms

- Logic simulation
- Packet processing
- Digital signal processing

■ Parallel microarchitecture

- A CPU+GPU integrated packet processing engine

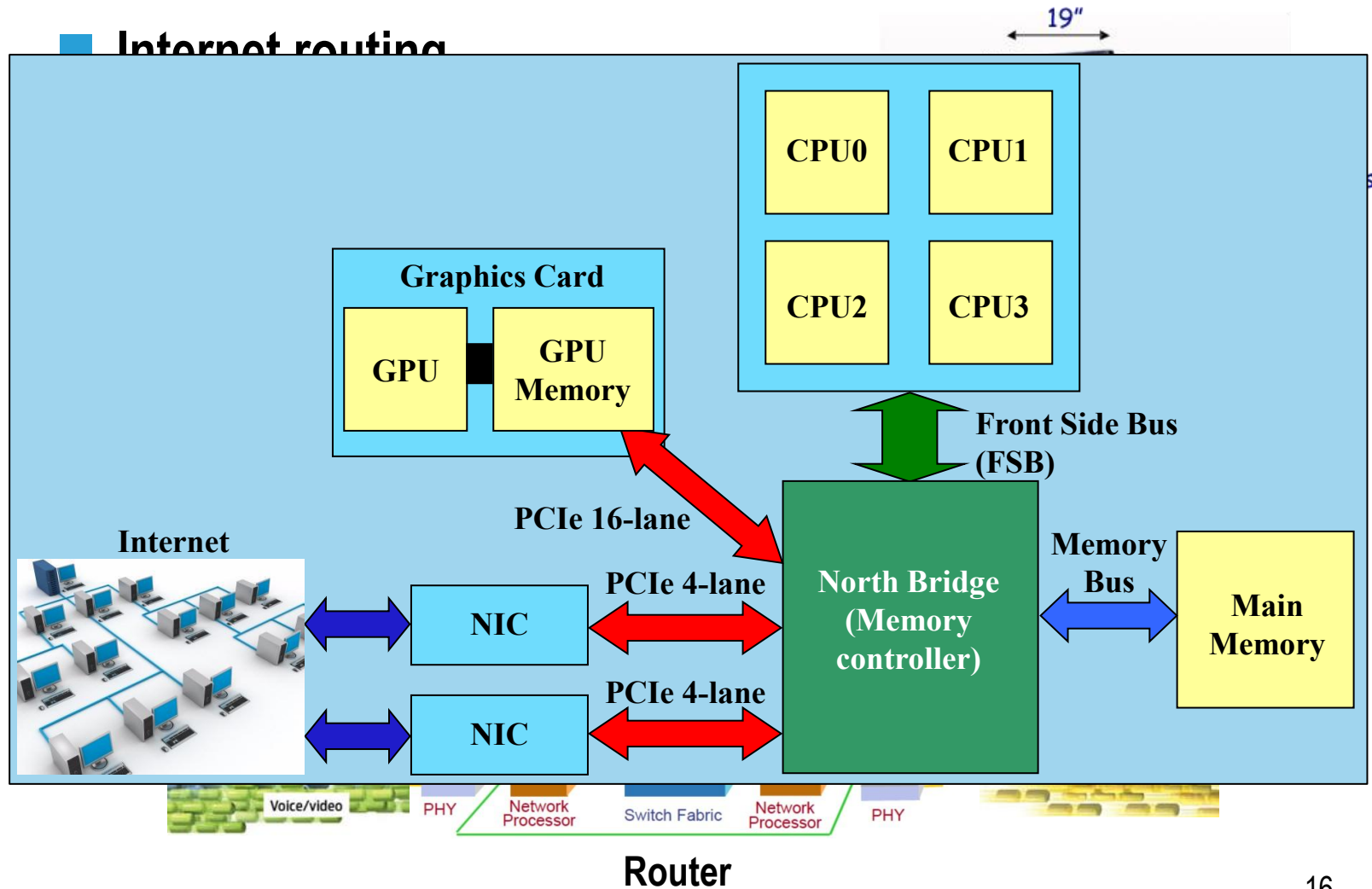
■ GPU programming

- Source-to-source GPU code generation

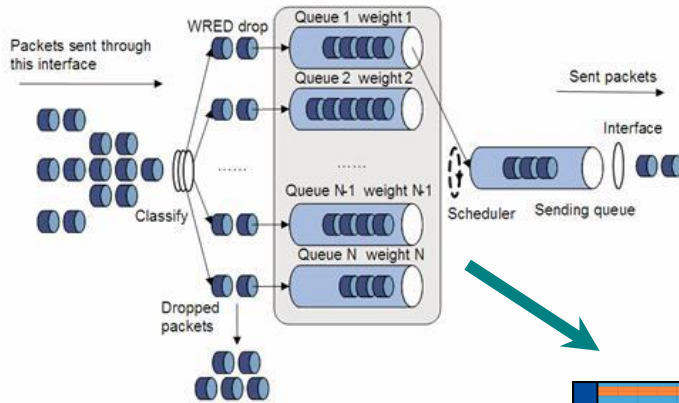


GPU Accelerated Software Router

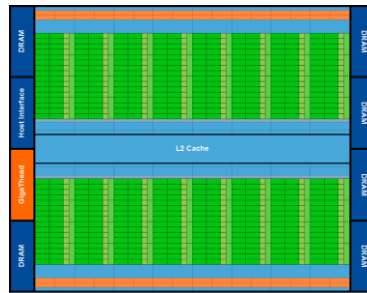
Internet routing



GPU Based Software Router



**Packet classification
(30X Speed-up)**



GPU

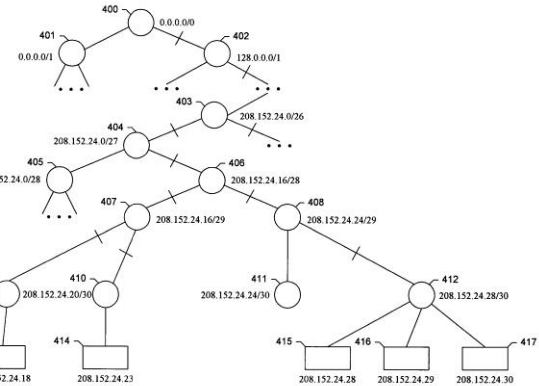
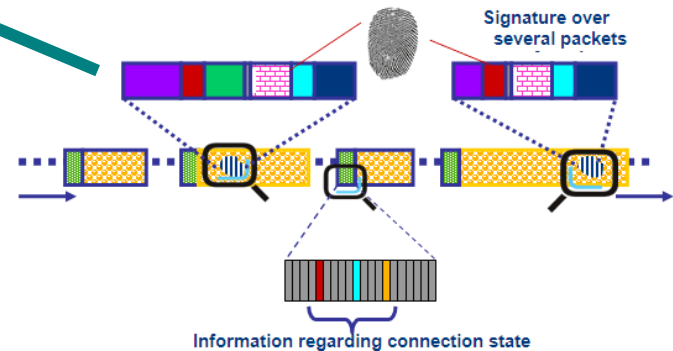


Table lookup (7X Speed-up)



**Network intrusion detection
(15-30X Speed-up)**

***Published on Design Automation & Test Europe
2010 and 2011**

Research Outline

■ Applications and algorithms

- Logic simulation
- Packet processing
- Digital signal processing

■ Parallel microarchitecture

- CPU+GPU integrated packet processing microarchitecture

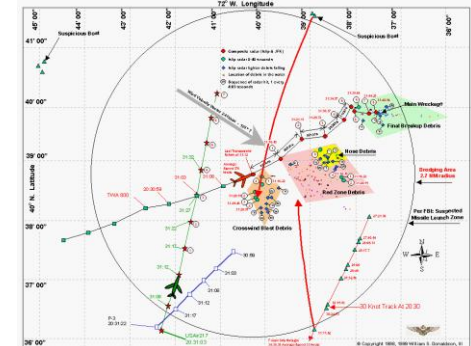
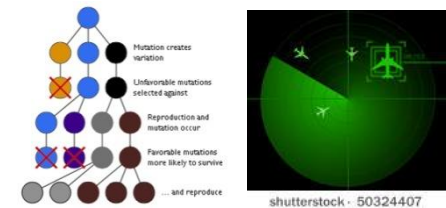
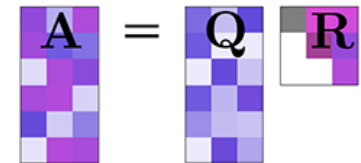
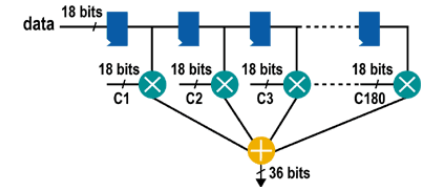
■ GPU programming

- Source-to-source GPU code generation



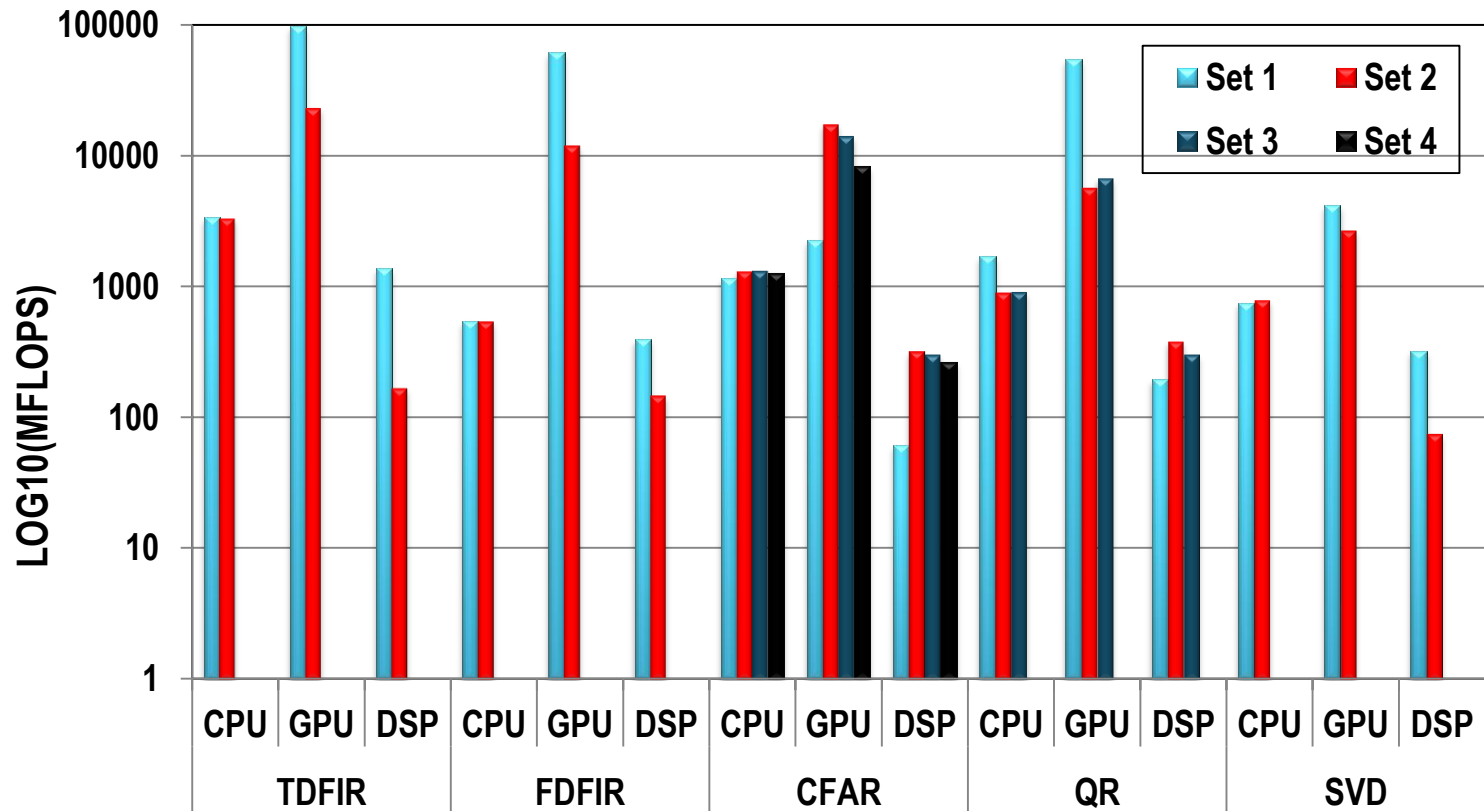
HPEC Challenge - Radar Benchmarks

Benchmark	Description
TDFIR	Time-domain finite impulse response filtering
FDFIR	Frequency-domain finite impulse response filtering
CT	Corner turn or matrix transpose to place radar data into a contiguous row for efficient FFT
QR	QR factorization: prevalent in target recognition algorithms
SVD	Singular value decomposition: produces a basis for the matrix as well as the rank for reducing interference
CFAR	Constant false-alarm rate detection: find target in an environment with varying background noise
GA	Graph optimization via genetic algorithm: removing uncorrelated data relations
PM	Pattern Matching: identify stored tracks that match a target
DB	Database operations to store and query target tracks



Performance Comparison

- GPU: NVIDIA Fermi, CPU: Intel Core 2 Duo (3.33GHz), DSP AD TigerSharc 101

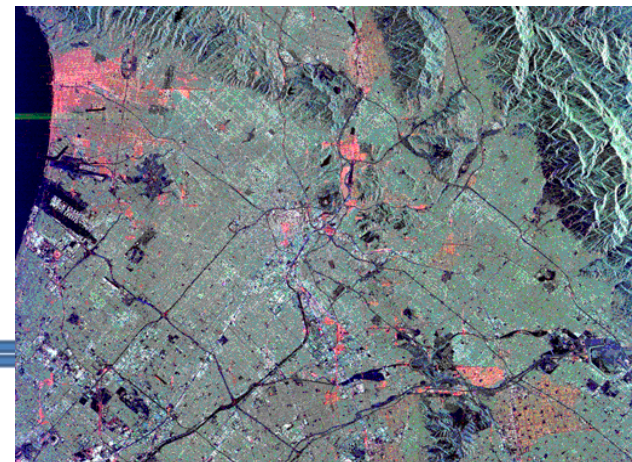
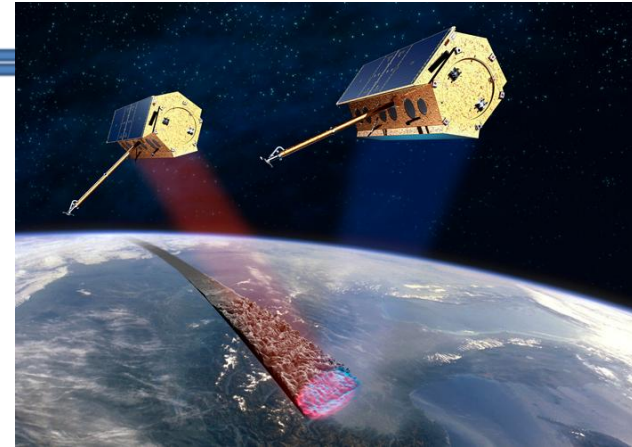
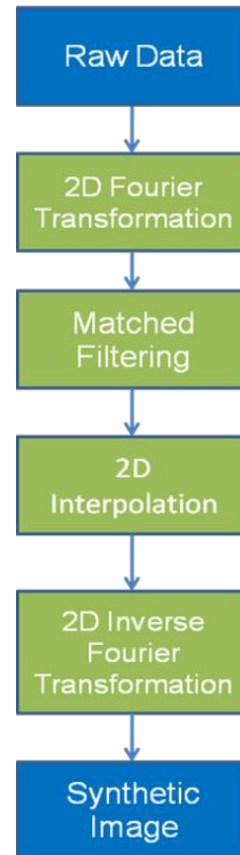


***Published on Design Automation & Test Europe 2011**

Synthetic Aperture Radar (SAR) on GPU

■ A complete radar application implemented on GPU

- Imaging radar
- ~30X speedup
- Performance results including CPU-GPU data transfer



Research Outline

■ Applications and algorithms

- Logic simulation
- Packet processing
- Digital signal processing

■ Parallel microarchitecture

- A CPU+GPU integrated packet processing engine

■ GPU programming

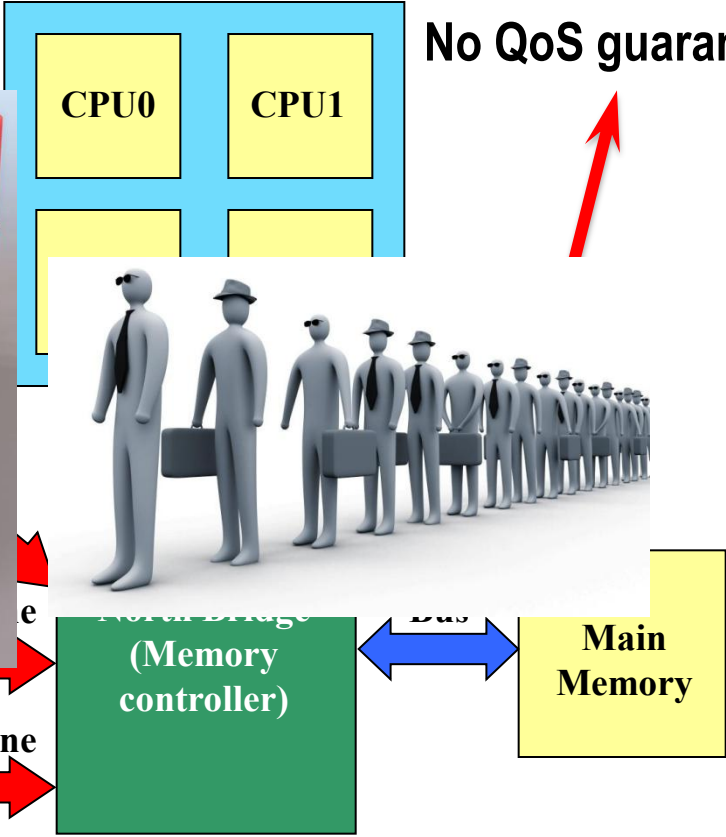
- Source-to-source GPU code generation



Limitation of GPU-Based Packet Processing

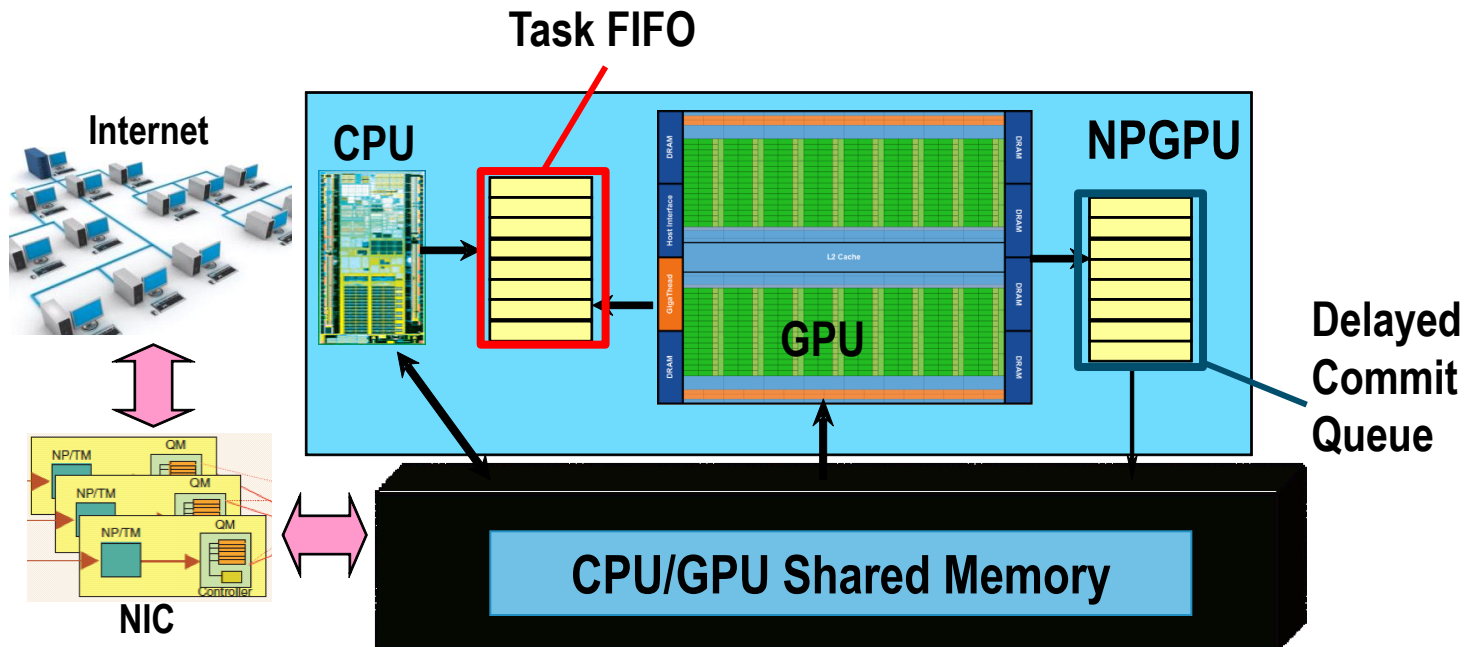
CPU-GPU communication overhead

No QoS guarantee



Morphing GPU into a Network Processor

- CPU-GPU integration
 - Shared memory space - 5X performance improvement
- Latency aware scheduling
 - Reduce packet latency by 82%



*Published on Design Automation Conference 2011

Research Outline

■ Applications and algorithms

- Logic simulation
- Packet processing
- Digital signal processing

■ Parallel microarchitecture

- A CPU+GPU integrated packet processing engine

■ GPU programming

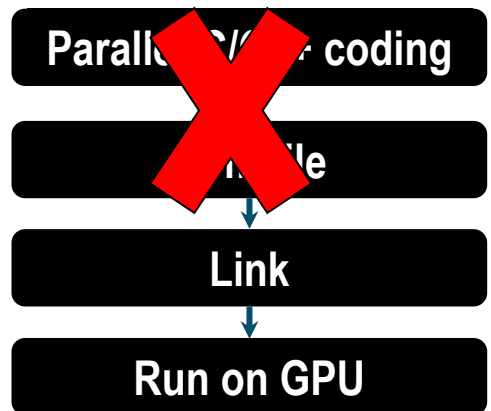
- Source-to-source GPU code generation



Source-to-Source GPU Code Generation

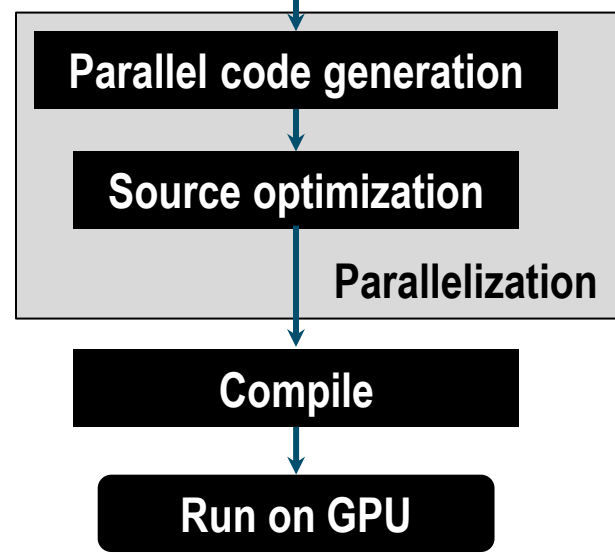
- GPU programming is challenging
 - Load balance, synchronization, hardware details
- Source level GPU code generation
 - Applications domains: Numerical, DSP, parallel embedded applications
 - Input: Algorithmic specific
 - Output: Parallel code on GPU

```
for (i=1; i<4; ++i)
  for (j=1; j<5; ++j) {
    A[i,j] = A[i-1,j-1];
  }
```



Current GPU programming flow

Algorithm specification or legacy code



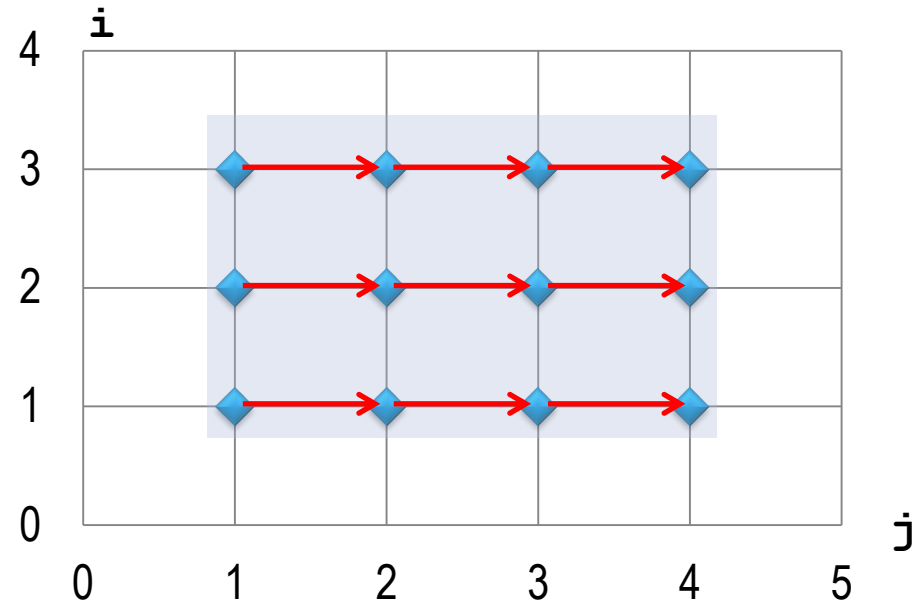
Our new GPU code generation flow

Automatic Parallelization

■ Polyhedron based loop parallelization

- Loop bounds defines a domain polyhedron
 - Discrete space
- Loop dependency represented by directed edges

```
for (i=1; i<=3; ++i)
  for (j=1; j<=4; ++j) {
    A[i,j] = A[i,j-1];
  }
```



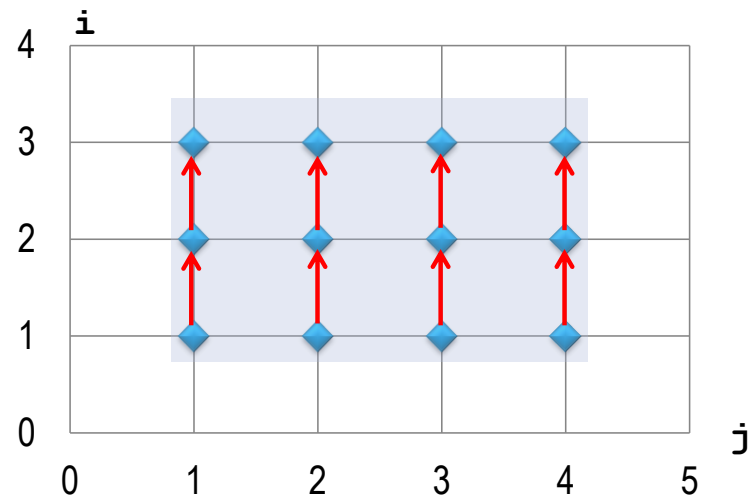
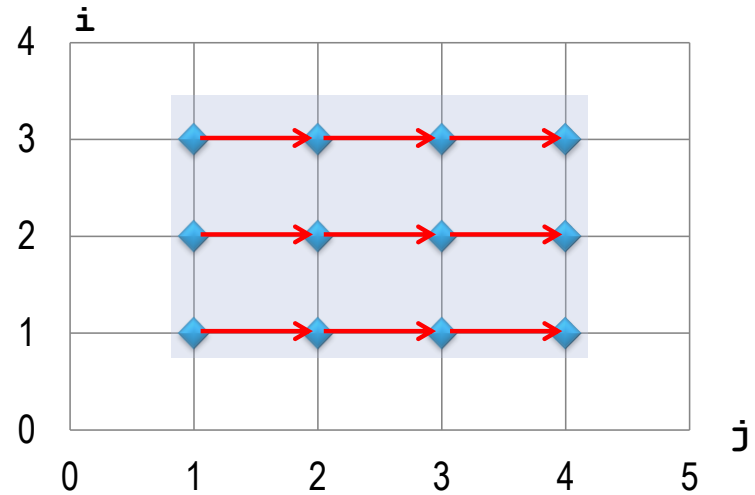
Polyhedron Based Loop Parallelization

```
for (i=1; i<=3; ++i)
  for (j=1; j<=4; ++j) {
    A[i,j] = A[i,j-1];
  }
```

Permutation

$$\begin{bmatrix} p \\ q \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} i \\ j \end{bmatrix}$$

```
for (p=1; p<=4; ++p)
  for (q=1; q<=3; ++q) {
    A[p,1] = A[i-1,j];
  }
```



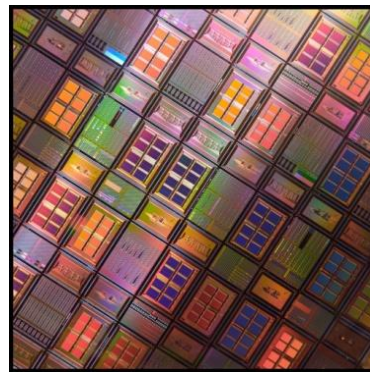
Outline



清华大学

Tsinghua University

1. Introduction



2. Research



3. CUDA Propagation

Related Classes

- **Prof. Yangdong Deng (Institute of Microelectronics)**
 - 5-day short courses
 - Offered at Tsinghua University and China Academy of Science (Supercomputing center and institute of Acoustics)
 - 400 attendees
- **Prof. Wei Xue & Yongwei Wu (CS Department)**
 - CUDA module in “Parallel Programming” (undergraduate) and “Parallel Programming Labs” (graduates)
- **Developing “GPU based Parallel Programming”**
 - Selected by 国家精品课程中心 (National Center of Excellent Courses)
 - Will be offered by major mainland China universities

Textbooks

- **Y. Deng and W. Liu, “Massively Data Parallel Algorithms,” Tsinghua Publishing House, in press.**
- **Y. Deng, Y. Liu, and H. Chen, “GPU Based Parallel Computing,” China Advanced Education Publisher, Textbook Series on Advanced Industry Technology Courses, in press.**

CUDA University Programming Contest

- Annual contest started in 2009
- Continued in 2010
 - 1,500+ students from 200+ schools registered
 - 122 programs submitted
 - Designated topics: 22
 - Self-proposed topics: 100
 - > A wide spectrum of application domains covered
 - Scientific/engineering/consumer applications



- 2011 contest incoming!

